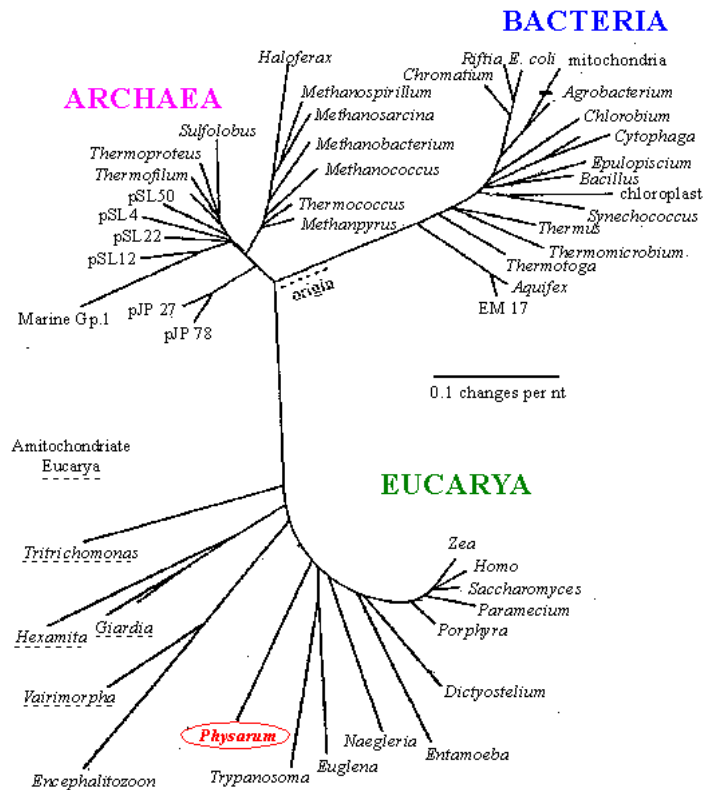

Pairwise Alignment

Anders Gorm Pedersen
Molecular Evolution Group
Center for Biological Sequence Analysis

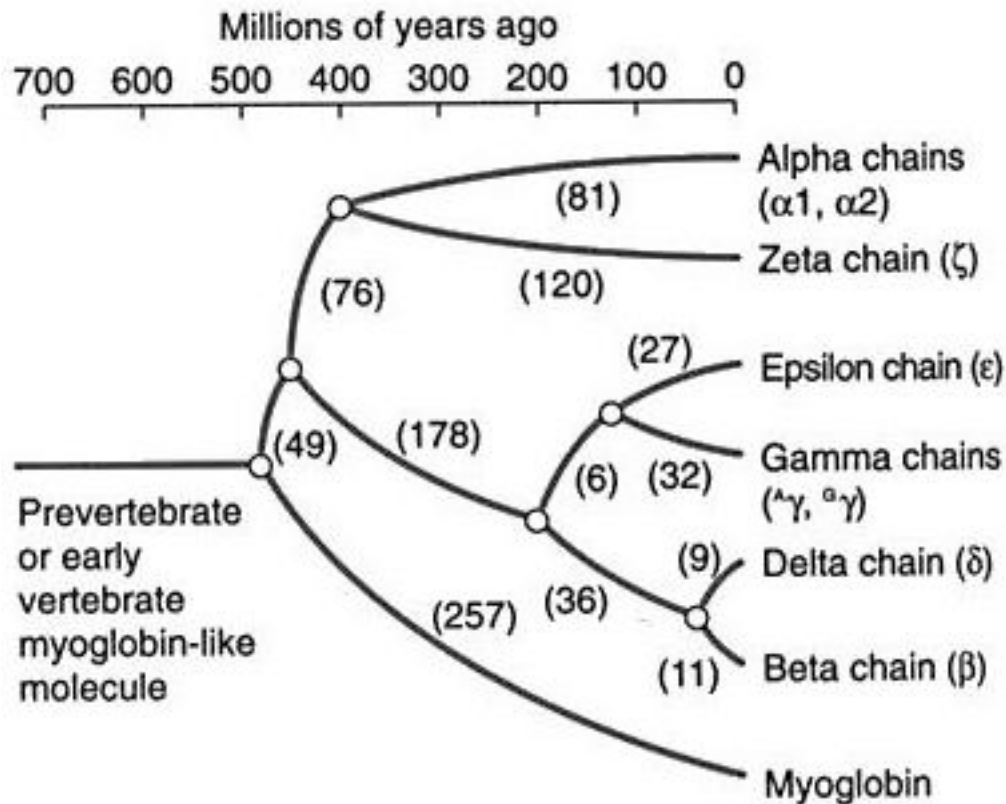
Sequences are related

- Darwin: all organisms are related through descent with modification
- => Sequences are related through descent with modification
- => Similar molecules have similar functions in different organisms



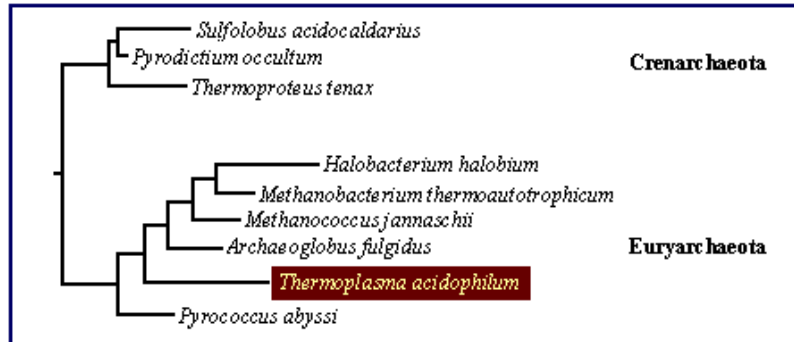
Phylogenetic tree based on
ribosomal RNA:
three domains of life

Sequences are related, II



Phylogenetic tree of globin-type proteins found in humans

Why compare sequences?



- Determination of evolutionary relationships

Protein 1: binds oxygen

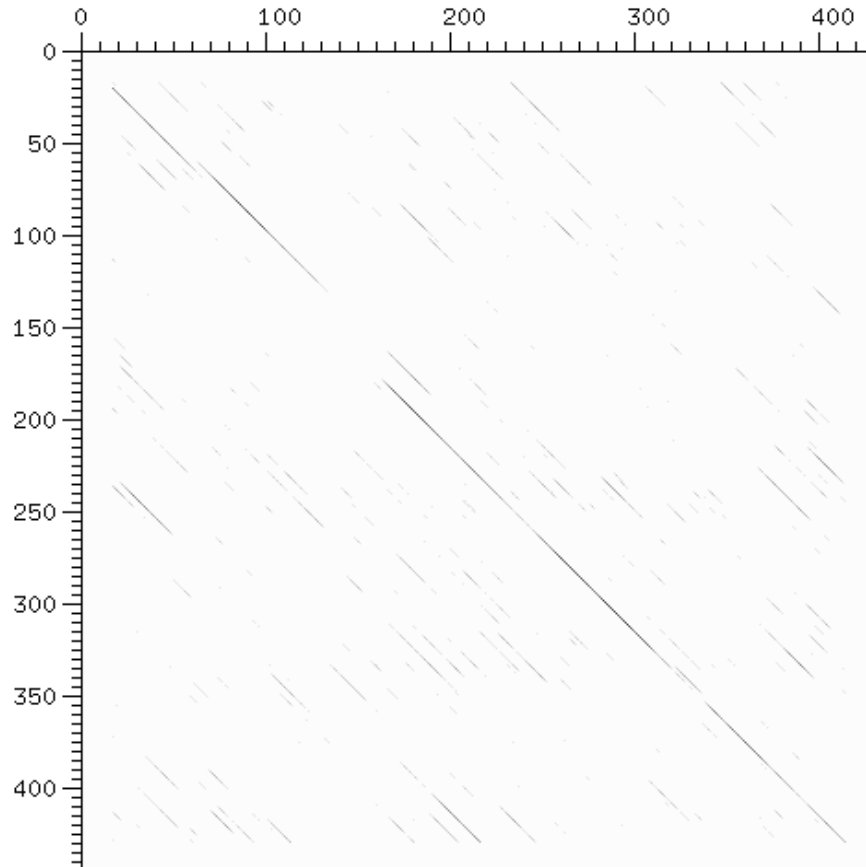


Sequence similarity

Protein 2: binds oxygen ?

- Prediction of protein function and structure (database searches).

Dotplots: visual sequence comparison



1. Place two sequences along axes of plot
2. Place dot at grid points where two sequences have identical residues
3. Diagonals correspond to conserved regions

Pairwise alignments

43.2% identity;

Global alignment score: 374

	10	20	30	40	50	
alpha	V-LSPADKTNVKA	AWGKVGAHAGEYGA	EALERMFLSFPTTK	TYFPHF-DLS----	HGSA	
	:	:::	...	:	:	:
beta	VHLTPEEKSAVTAL	WGKV--NVDEVGGEAL	GRLLVVYPWTQRFF	ESFGDLSTPD	AVMGNP	
	10	20	30	40	50	
	60	70	80	90	100	110
alpha	QVKGHGKKVADALT	NAVAHVDDMPNALS	SALSDLHAHKLRVDP	VNFKLLSHCLLVTL	AHL	

beta	KVKAHGKKVLGAFSD	GLAHLNLTGTFATL	SELHCDKLHVDPEN	FRLLGNVLCVLA	HHF	
	60	70	80	90	100	110
	120	130	140			
alpha	PAEFTPAVHASLDK	FLASVSTVLTSKYR				
	::::	
beta	GKEFTPPVQAAYQK	VVAGVANALAHKYH				
	120	130	140			

Pairwise alignment

100.000% identity in 3 aa overlap

SPA

:::

SPA

Percent identity is not a good measure of alignment quality

Global alignment score: 374

		10	20	30	40	50
alpha	V-LSPADKTNVKA	AWGKVG	AHAGEYGA	EALERMFLSF	PTTKTYFP	HF-DLS-----HGSA
	:	:::	..	:	:	:
beta	VH	LTPEEKSAVT	ALWGKV--	NVDEVGGE	ALGRLLV	VPWTQRFFESFGDLSTPDAVMGNP
		10	20	30	40	50
		60	70	80	90	100
alpha	QVK	GHGKKV	ADALTNA	VAHVDD	MPNALS	SALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHL
	:	:	:	:	:	:
beta	KV	KAHGKKV	LGAFSD	GLAHL	DNLKG	TFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHF
	60	70	80	90	100	110
		120	130	140		
alpha	PAE	FTPAV	HASLDK	FLASV	STVLT	SKYR
	:	:	:	:	:	:
beta	GKE	FTPPV	QAAYQ	KVVAG	VANAL	AHKYH
	120	130	140			

Alignment scores: match vs. mismatch

Simple scoring scheme (too simple in fact...):

Matching amino acids: 5

Mismatch: 0

Scoring example:

K A W S A D V

: : : : :

K D W S A E V

5+0+5+5+5+0+5 = 25

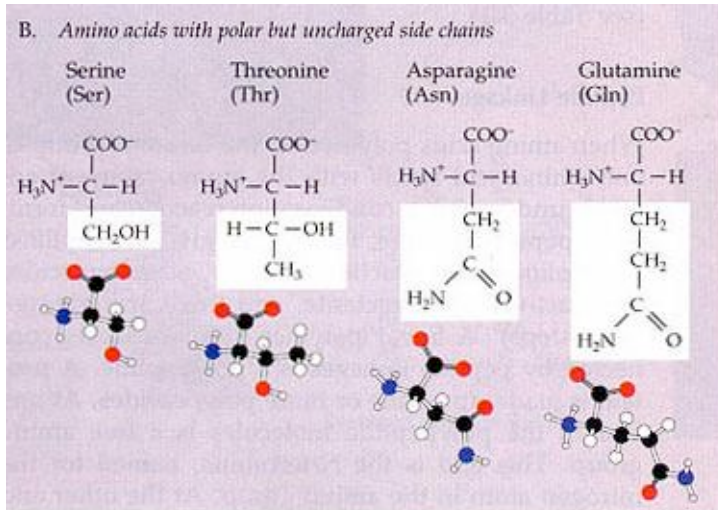
Pairwise alignments: conservative substitutions

43.2% identity;

Global alignment score: 374

	10	20	30	40	50			
alpha	V-LSPADKTNVKA	AWGKVG	AHAGEYGA	EALERMFLSF	PTTKTYFPHF-DLS-----HGSA			
	: : . : . : : : : : . : : : : : : : . : : : : : .							
beta	VHLTPEEKSAVT	ALWGKV--	NVDEVGGEAL	GRLLVVYPWT	QRFFESFGDLSTPD			
	10	20	30	40	50			
	60	70	80	90	100	110		
alpha	QVKGHGKKVAD	ALTNAVAHV	DDMPNALS	SALSDLHAH	KLRVDPVNF	KLLSHCLLV	TAAHL	
	: .							
beta	KVKAHGKKVL	GAFSDGLA	HLDNLKG	TFATLSEL	HCDKLH	VDPENFR	LLGNVLC	VLAHHF
	60	70	80	90	100	110		
	120	130	140					
alpha	PAEFTPAVHAS	LDKFLAS	VSTVLTSKYR					
	: : : : : : . : : : : : : : .							
beta	GKEFTPPVQ	AAYQKV	VAGVANAL	AHKYH				
	120	130	140					

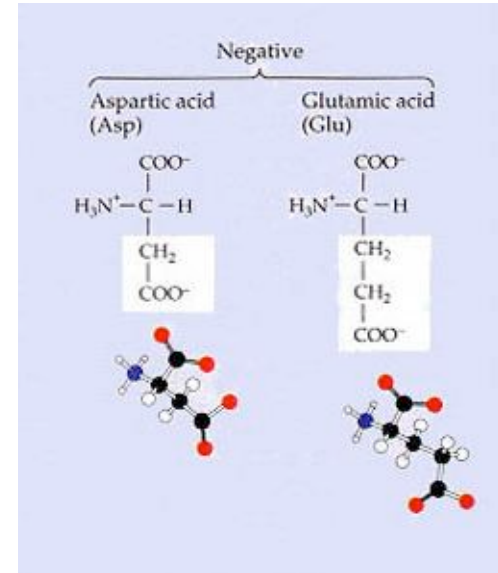
Amino acid properties



Serine (S) and Threonine (T) have similar physicochemical properties

=> Substitution of S/T or E/D occurs relatively often during evolution

=> Substitution of S/T or E/D should result in scores that are only moderately lower than identities



Aspartic acid (D) and Glutamic acid (E) have similar properties

Protein substitution matrices

A	5																			
R	-2	7																		
N	-1	-1	7																	
D	-2	-2	2	8																
C	-1	-4	-2	-4	13															
Q	-1	1	0	0	-3	7														
E	-1	0	0	2	-3	2	6													
G	0	-3	0	-1	-3	-2	-3	8												
H	-2	0	1	-1	-3	1	0	-2	10											
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5										
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5									
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6								
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7							
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8						
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10					
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5				
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5			
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-3	15		
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V

BLOSUM50 matrix:

- Positive scores on diagonal (identities)
- Similar residues get higher (positive) scores
- Dissimilar residues get smaller (negative) scores

Pairwise alignments: insertions/deletions

43.2% identity;

Global alignment score: 374

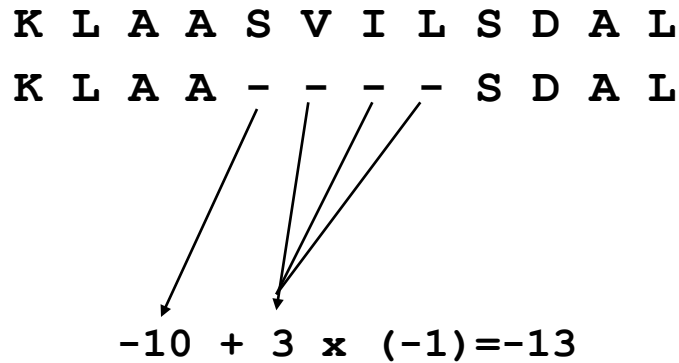
```

      10      20      30      40      50
alpha  V-LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLS-----HGSA
      :  ::  ::  :  :  ::  ..  :  ::  ::  ::  ::  :  :  :  :  :  :  :  :  :
beta   VHLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP
      10      20      30      40      50

      60      70      80      90     100     110
alpha  QVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHL
      .....  .....  .....  .....  .....  ..  ::  ::
beta   KVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLCVLAHHF
      60      70      80      90     100     110

      120     130     140
alpha  PAEFTPAVHASLDKFLASVSTVLTSKYR
      ::  ::  :  .....  ::
beta   GKEFTPPVQAAYQKVVAGVANALAHKYH
      120     130     140
```

Alignment scores: insertions/deletions



Affine gap penalties:

Multiple insertions/deletions may be one evolutionary event =>

Separate penalties for **gap opening** and **gap elongation**

Handout

Compute 4 alignment scores: two different alignments using two different alignment matrices (and the same gap penalty system)

Score 1: Alignment 1 + BLOSUM-50 matrix + gaps

Score 2: Alignment 1 + BLOSUM-Trp matrix + gaps

Score 3: Alignment 2 + BLOSUM-50 matrix + gaps

Score 4: Alignment 2 + BLOSUM-Trp matrix + gaps



Note: fake matrix constructed for pedagogic purposes.

Handout: summary of results

	Alignment 1	Alignment 2
BLOSUM-50	38	51
BLOSUM-Trp	118	91

Protein substitution matrices

A	5																			
R	-2	7																		
N	-1	-1	7																	
D	-2	-2	2	8																
C	-1	-4	-2	-4	13															
Q	-1	1	0	0	-3	7														
E	-1	0	0	2	-3	2	6													
G	0	-3	0	-1	-3	-2	-3	8												
H	-2	0	1	-1	-3	1	0	-2	10											
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5										
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5									
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6								
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7							
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8						
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10					
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5				
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5			
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-3	15		
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V

BLOSUM50 matrix:

- Positive scores on diagonal (identities)
- Similar residues get higher (positive) scores
- Dissimilar residues get smaller (negative) scores

Protein substitution matrices: different types

- **Identity matrix**
(match vs. mismatch)
- **Chemical properties matrix**
(use knowledge of physicochemical properties to design matrix)
- ➔ • **Empirical matrices**
(based on observed pair-frequencies in hand-made alignments)
 - PAM series
 - BLOSUM series
 - Gonnet

Estimation of the BLOSUM 50 matrix

- BLOSUM matrices are computed based on gap-free alignments in the so-called BLOCKS database. BLOSUM 50 is computed by comparing sequences that are less than 50% identical. BLOSUM 80 is computed from sequences less than 80% identical, etc.
- All pairs of sequences in a block are compared, and the **observed pair frequencies** (p_{ab}) are noted. For instance: $p_{WW} = 0.0065$, $p_{AL} = 0.0044$, etc.
- Expected pair frequencies** are computed from single amino acid frequencies. For instance: $p_A \times p_L = 0.074 \times 0.099 = 0.0073$
- For each amino acid pair the substitution scores are essentially computed as follows (here λ is a scaling factor, used to obtain integer scores):

$$S(a, b) = \frac{1}{\lambda} \ln \left(\frac{p_{\text{obs}}}{p_{\text{exp}}} \right) = \frac{1}{\lambda} \ln \left(\frac{p_{ab}}{p_a \times p_b} \right)$$

```
ID    FIBRONECTIN_2; BLOCK
COG9_CANFA  GNSAGEPCVFPFIFLGKQYSTCTREGRGDGHLWCATT
COG9_RABIT  GNADGAPCHFPFTFEGRSYTACTTDGRSDGMAWCSTT
FA12_HUMAN  LTVTGEPCHFPPFQYHRQLYHKCTHKGRPGPQWCATT
HGFA_HUMAN  LTEDGRPCRFPPFRYGGRLHACTSEGAHRKWCATTH
MANR_HUMAN  GNANGATCAFPFKFENKQYADCTSGRSDGWLWCGTT
MPRI_MOUSE  ETDDGEPVFPFPIYKGSYDECVLGRAKLWCSKTAN
PB1_PIG     AITSDDKCVFPFIYKGNLYFDCTLHDSTYYWCSVTY
SFP1_BOVIN  ELPEDEECVFPFVYRNRKHFDCVHGSFLFPWCSLDAD
SFP3_BOVIN  AETKDNKCVFPFIYGNKKYFDCTLHGSFLFLWCSLDAD
SFP4_BOVIN  AVFEGPACAFPTYKGGKYYMCTRKNSVLLWCSLDTE
SP1_HORSE   AATDYAKCAFPFVYRGQTYDRCTTDGSLFRISWCSVT
COG2_CHICK  GNSEGAPCVFPFIFLGNKYDSCTSGRNDGKLWCAST
COG2_HUMAN  GNSEGAPCVFPFTFLGNKYESCTSGRSDGKMWCAAT
COG2_MOUSE  GNSEGAPCVFPFTFLGNKYESCTSGRNDGKVWCATT
COG2_RABIT  GNSEGAPCVFPFTFLGNKYESCTSGRSDGKMWCATS
COG2_RAT    GNSEGAPCVFPFTFLGNKYESCTSGRNDGKVWCATT
COG9_BOVIN  GNADGKPCVFPFTFQGRYSACTSDGRSDGYRWCATT
COG9_HUMAN  GNADGKPCQFPFIFQGSYSACTTDGRSDGYRWCATT
COG9_MOUSE  GNGEGKPCVFPFIFEGRSYSACTTKGRSDGYRWCATT
COG9_RAT    GNGDGKPCVFPFIFEGHSYSACTTKGRSDGYRWCATT
FINC_BOVIN  GNSNGALCHFPFLYNNHNYTDCTSEGRDNDNMKWCATT
FINC_HUMAN  GNSNGALCHFPFLYNNHNYTDCTSEGRDNDNMKWCATT
FINC_RAT    GNSNGALCHFPFLYNNHNYSDCTSEGRDNDNMKWCATT
MPRI_BOVIN  ETEDGEPCVFPFVFNKGSYEECVVESRARLWCATTAN
MPRI_HUMAN  ETDDGVPCVFPFIFNKGSYEECIIESRAKLWCSTTAD
PA2R_BOVIN  GNAHGTPCMFPFQYNQQWHHECTREGREDNLLWCATT
PA2R_RABIT  GNAHGTPCMFPFQYNHQQWHHECTREGRQDDSLWCATT
```

Estimation of the BLOSUM 50 matrix

- Example: A + L:

$$p_A = 0.074, p_L = 0.099 \Rightarrow p_A \times p_L = 0.0073$$

$$p_{AL} = 0.044$$

$$\begin{aligned} S_{AL} &= \frac{1}{\lambda} \ln \left(\frac{p_{AL}}{p_A \times p_L} \right) \\ &= \frac{1}{0.347} \ln \left(\frac{0.0044}{0.0073} \right) \\ &= \frac{-0.51}{0.347} \\ &= -1.46 \\ &\approx -1 \end{aligned}$$

Pairwise alignment

Optimal alignment:

alignment having the highest possible score given a substitution matrix and a set of gap penalties

So:

best alignment can be found by exhaustively searching all possible alignments, scoring each of them and choosing the one with the highest score?

How many possible alignments are there?

Handout exercise: Enumerate all the possible alignments for two sequences of lengths $n_1=2$ and $n_2=3$

How many possible alignments are there?

Derivation of formula

- **Starting point:** Two sequences, s_1 and s_2
 - Lengths: n_1 and n_2 , where s_1 is the longest ($n_1 > n_2$)

How many possible alignments are there?

Derivation of formula

- **Starting point:** Two sequences, s_1 and s_2
 - Lengths: n_1 and n_2 , where s_1 is the longest ($n_1 > n_2$)

- Example:

s_1 : ABCDE $n_1=5$

s_2 : 123 $n_2=3$

How many possible alignments are there?

Derivation of formula

- **Step 1:** We want to add k gaps to s_1
 - Note: $k \leq n_2$ (we don't want to align gap with gap)
 - The total length of the alignment will then be: $n_1 + k$

How many possible alignments are there?

Derivation of formula

- **Step 1:** We want to add k gaps to s_1
 - Note: $k \leq n_2$ (we don't want to align gap with gap)
 - The total length of the alignment will then be: $n_1 + k$

$$k=2 \Rightarrow n_1 + k = 7$$

(potential alignment positions shown as empty boxes below)

s_1 : ABCDE

s_2 : 123

How many possible alignments are there?

Derivation of formula

- **Step 2:** Place the n_2 non-gap symbols of s_2 in the n_1+k boxes
 - The number of ways of doing this is given by the following binomial coefficient:

$$N = \binom{n_1 + k}{n_2} = \frac{(n_1 + k)!}{(n_2)!(n_1 + k - n_2)!}$$

s_1 : ABCDE

s_2 : 123

How many possible alignments are there?

Derivation of formula

- **Step 2:** Place the n_2 non-gap symbols of s_2 in the boxes
 - The number of ways of doing this is given by the following binomial coefficient:

$$N = \binom{n_1 + k}{n_2} = \frac{(n_1 + k)!}{(n_2)!(n_1 + k - n_2)!}$$

s_1 : ABCDE

s_2 :

		1		2	3	

$$N = \frac{7!}{3! \times 4!} = 35$$

How many possible alignments are there?

Derivation of formula

- **Step 2:** Place the n_2 non-gap symbols of s_2 in the boxes
 - The number of ways of doing this is given by the following binomial coefficient:

$$N = \binom{n_1 + k}{n_2} = \frac{(n_1 + k)!}{(n_2)!(n_1 + k - n_2)!}$$

s_1 : ABCDE

s_2 :

-	-	1	-	2	3	-

$$N = \frac{7!}{3! \times 4!} = 35$$

- The remaining boxes will then contain gap symbols

How many possible alignments are there?

Derivation of formula

- **Step 3:** Place the k gaps in s_1 opposite the n_2 non-gap symbols in s_2 (this ensures that gaps in s_1 are never aligned with a gap in s_2)
- The number of ways of doing this is given by the following binomial coefficient:

$$N = \binom{n_2}{k} = \frac{(n_2)!}{k!(n_2 - k)!}$$

s_1 : ABCDE

s_2 :

-	-	1	-	2	3	-

How many possible alignments are there?

Derivation of formula

- **Step 3:** Place the k gaps in s_1 opposite the n_2 non-gap symbols in s_2 (this ensures that a gap in s_1 never is aligned with a gap in s_2)
- The number of ways of doing this is given by the following binomial coefficient:

$$N = \binom{n_2}{k} = \frac{(n_2)!}{k!(n_2 - k)!}$$

s_1 : ABCDE

s_2 :

		-			-	
-	-	1	-	2	3	-

$$N = \binom{3}{2} = \frac{3!}{2!1!} = 3$$

How many possible alignments are there?

Derivation of formula

- **Step 3:** Place the k gaps in s_1 opposite the n_2 non-gap symbols in s_2 (this ensures that gaps in s_1 are never aligned with a gap in s_2)
- The number of ways of doing this is given by the following binomial coefficient:

$$N = \binom{n_2}{k} = \frac{(n_2)!}{k!(n_2 - k)!}$$

s_1 :	A	B	-	C	D	-	E
s_2 :	-	-	1	-	2	3	-

$$N = \binom{3}{2} = \frac{3!}{2!1!} = 3$$

- The remaining boxes will then contain the non-gap symbols, in order

How many possible alignments are there?

Derivation of formula

- For k gaps added to s_1 we therefore have:
 - Number of ways of placing the n_2 non-gap symbols from s_2 in the n_1+k “alignment boxes”:

$$N = \binom{n_1 + k}{n_2}$$

-	-	1	-	2	3	-

- Number of ways of placing the k gap symbols from s_1 opposite from the n_2 non-gap symbols in this alignment:

$$N = \binom{n_2}{k}$$

A	B	-	C	D	-	E
-	-	1	-	2	3	-

- The total number of alignments with k gaps added to s_1 is therefore:

$$N = \binom{n_1 + k}{n_2} \binom{n_2}{k}$$

How many possible alignments are there?

Derivation of formula

- For each possible value of k the number of alignments is given by:

$$N = \binom{n_1 + k}{n_2} \binom{n_2}{k}$$

- Given that k can be any number from 0 to n_2 , the total number of possible alignments of sequences of length n_1 and n_2 is therefore:

$$N = \sum_{k=0}^{n_2} \binom{n_1 + k}{n_2} \binom{n_2}{k}$$

How many possible alignments are there?

Derivation of formula

- Example: $n_1=3$, $n_2=2$:

$$\begin{aligned} N &= \sum_{k=0}^2 \binom{3+k}{2} \binom{2}{k} \\ &= \binom{3}{2} \binom{2}{0} + \binom{4}{2} \binom{2}{1} + \binom{5}{2} \binom{2}{2} \\ &= 3 \times 1 + 6 \times 2 + 10 \times 1 \\ &= 3 + 12 + 10 \\ &= 25 \end{aligned}$$

How many possible alignments are there?

Length of sequences: $n_1 = n_2$	Number of possible alignments
2	13
3	63
4	321
5	1683
10	8,097,453
20	2.61×10^{14}
100	2.05×10^{75}
300	1.53×10^{228}

Pairwise alignment: the problem

The number of possible pairwise alignments increases explosively with the length of the sequences:

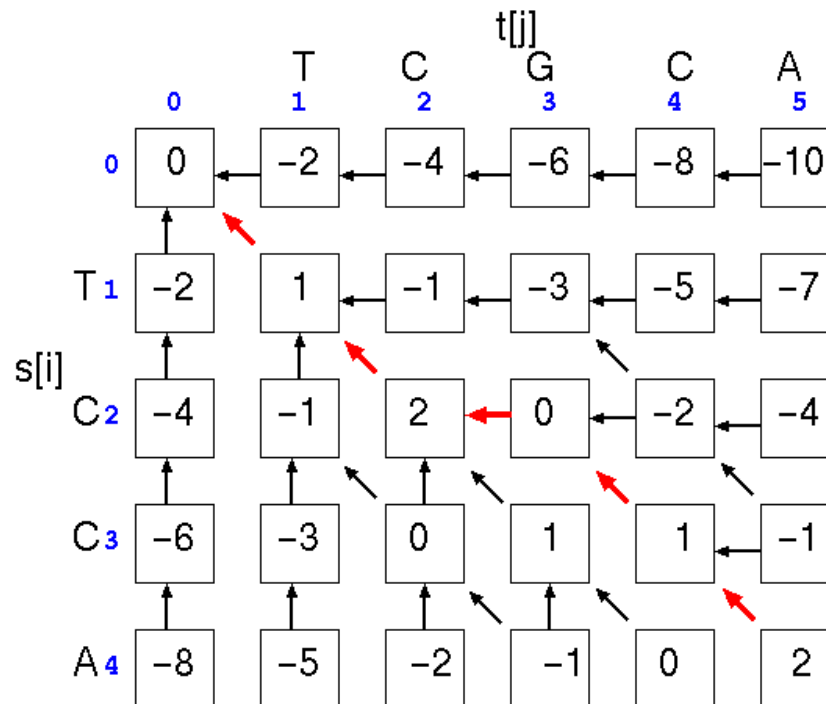
$$N = \sum_{k=0}^{n_2} \binom{n_1 + k}{n_2} \binom{n_2}{k}$$

Two protein sequences of length 300 amino acids can be aligned in approximately 10^{228} different ways

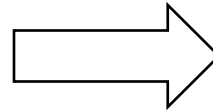
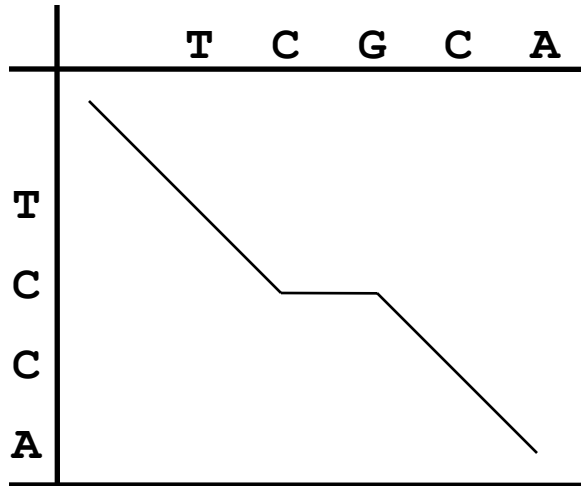
Time needed to test all possibilities much larger than the entire lifetime of the universe.

Pairwise alignment: the solution

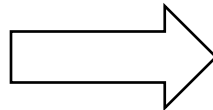
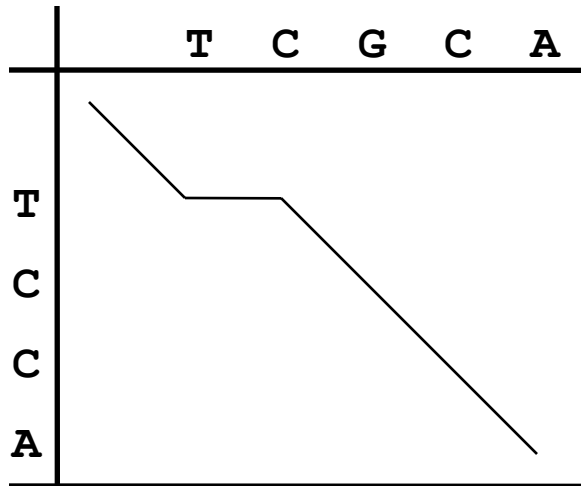
"Dynamic programming"
(the Needleman-Wunsch algorithm)



Alignment depicted as path in matrix



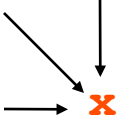
TCGCA
TC-CA



TCGCA
T-CCA

Dynamic programming: computation of scores

	T	C	G	C	A
T					
C					
C					
A					



The diagram shows a 5x5 matrix with columns labeled T, C, G, C, A and rows labeled T, C, C, A. The cell at row 2, column 2 (C, C) is marked with a red 'x'. Three arrows point to this cell: one from the cell above (T, C), one from the cell to the left (C, T), and one from the cell diagonally above and to the left (T, T).

Any given point in matrix can only be reached from three possible previous positions (you cannot “align backwards”).

=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

Dynamic programming: computation of scores

	T	C	G	C	A
T					
C		x			
C					
A					

Any given point in matrix can only be reached from three possible positions (you cannot “align backwards”).

=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

$$\text{score}(x,y) = \max \left\{ \begin{array}{l} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y) \\ \text{score}(x-1,y-1) \end{array} \right.$$

Dynamic programming: computation of scores

	T	C	G	C	A
T					
C					
C					
A					

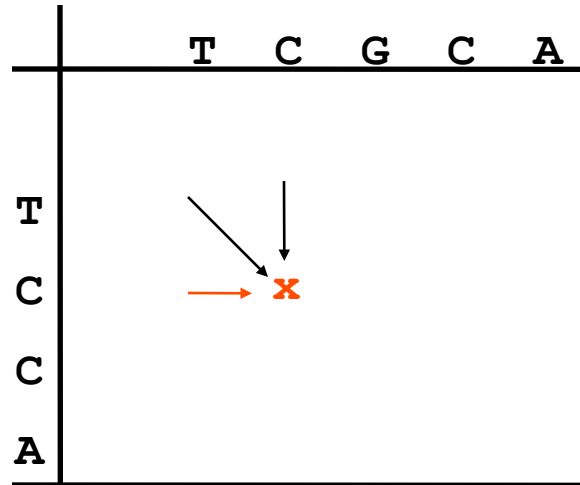
Any given point in matrix can only be reached from three possible positions (you cannot “align backwards”).

=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

$$\text{score}(x,y) = \max \left\{ \begin{array}{l} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \end{array} \right.$$

Dynamic programming: computation of scores

	T	C	G	C	A
T					
C					
C					
A					



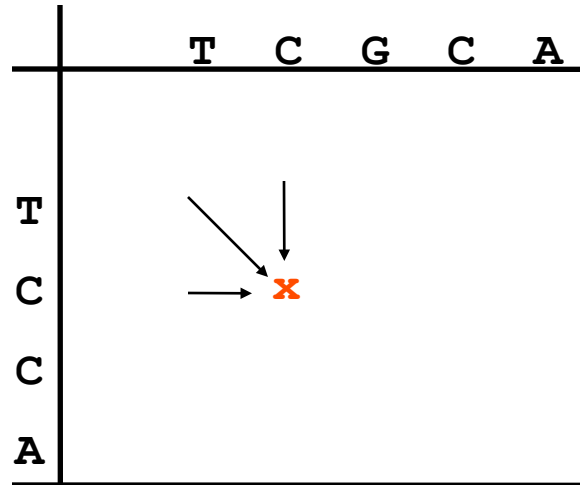
Any given point in matrix can only be reached from three possible positions (you cannot “align backwards”).

=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \\ \text{score}(x-1,y) - \text{gap-penalty} \end{cases}$$

Dynamic programming: computation of scores

	T	C	G	C	A
T					
C					
C					
A					



Any given point in matrix can only be reached from three possible positions (you cannot “align backwards”).

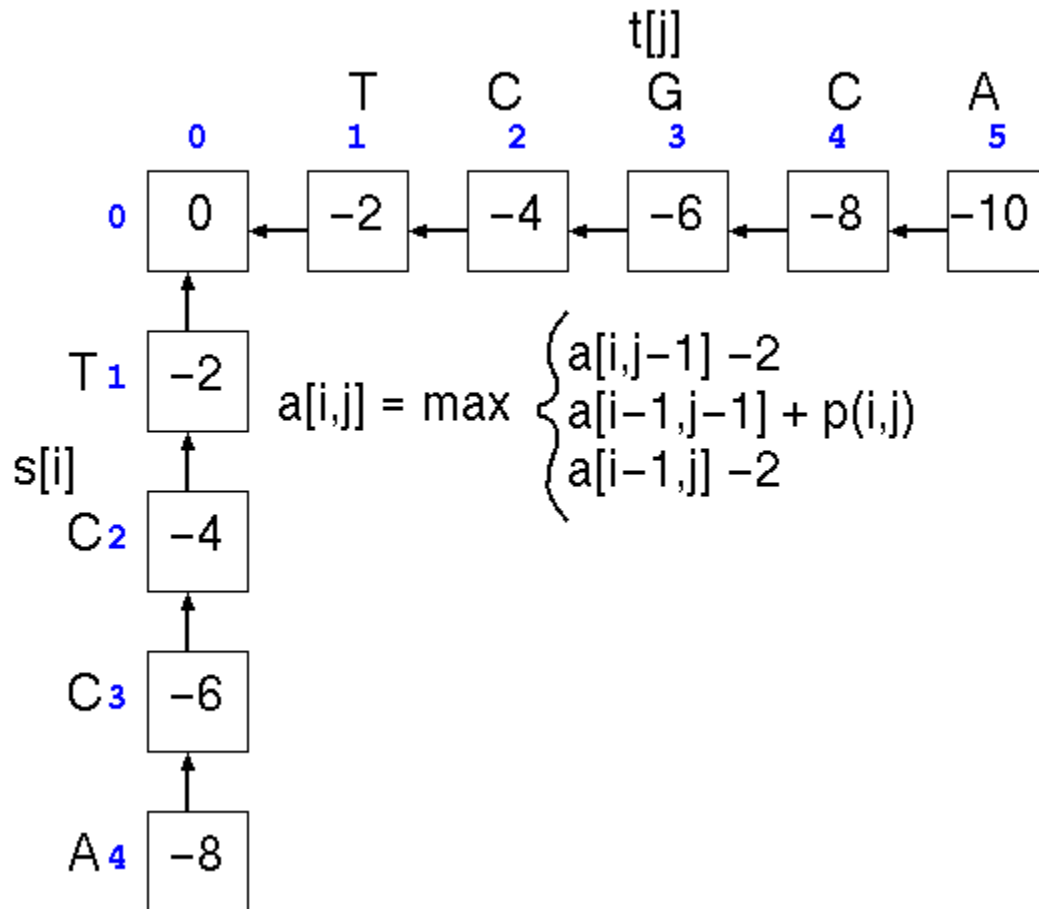
=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

Each new score is found by choosing the maximum of three possibilities.
For each square in matrix: keep track of where best score came from.

Fill in scores one row at a time, starting in upper left corner of matrix, ending in lower right corner.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \\ \text{score}(x-1,y) - \text{gap-penalty} \end{cases}$$

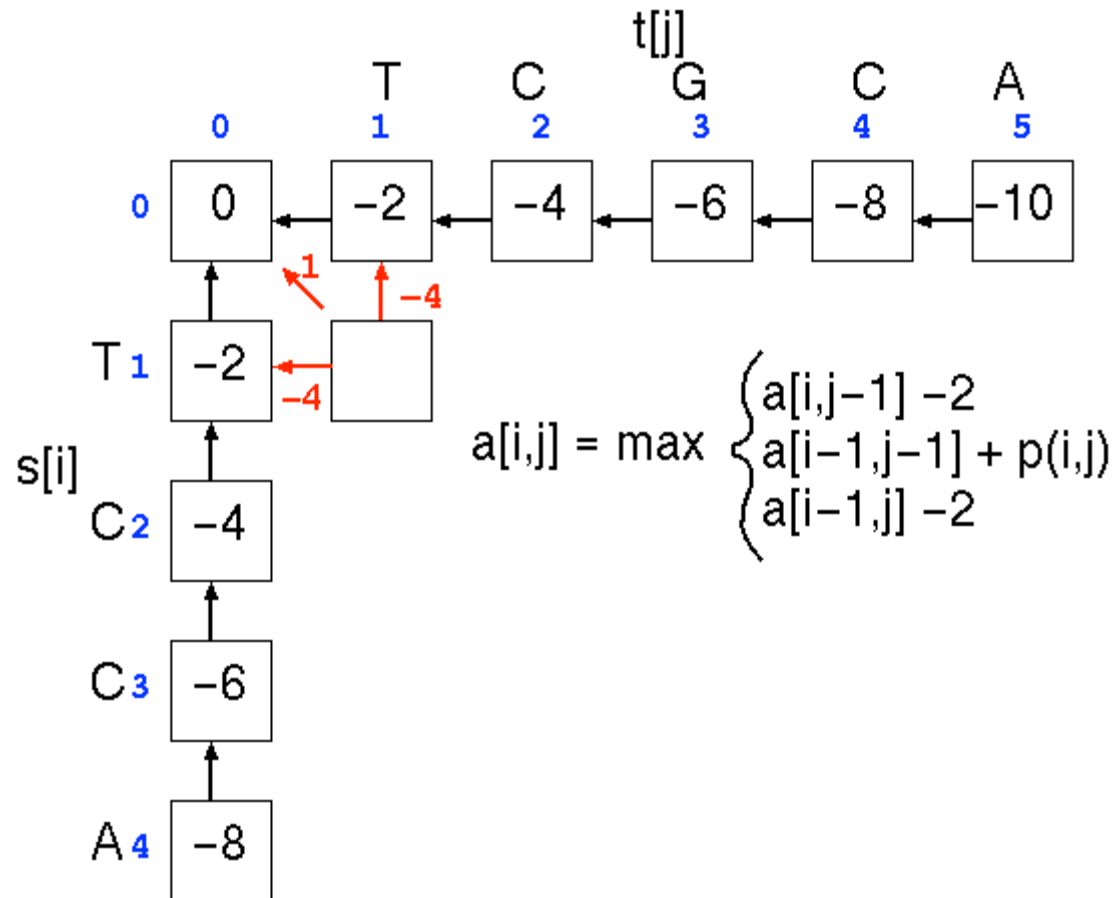
Dynamic programming: example



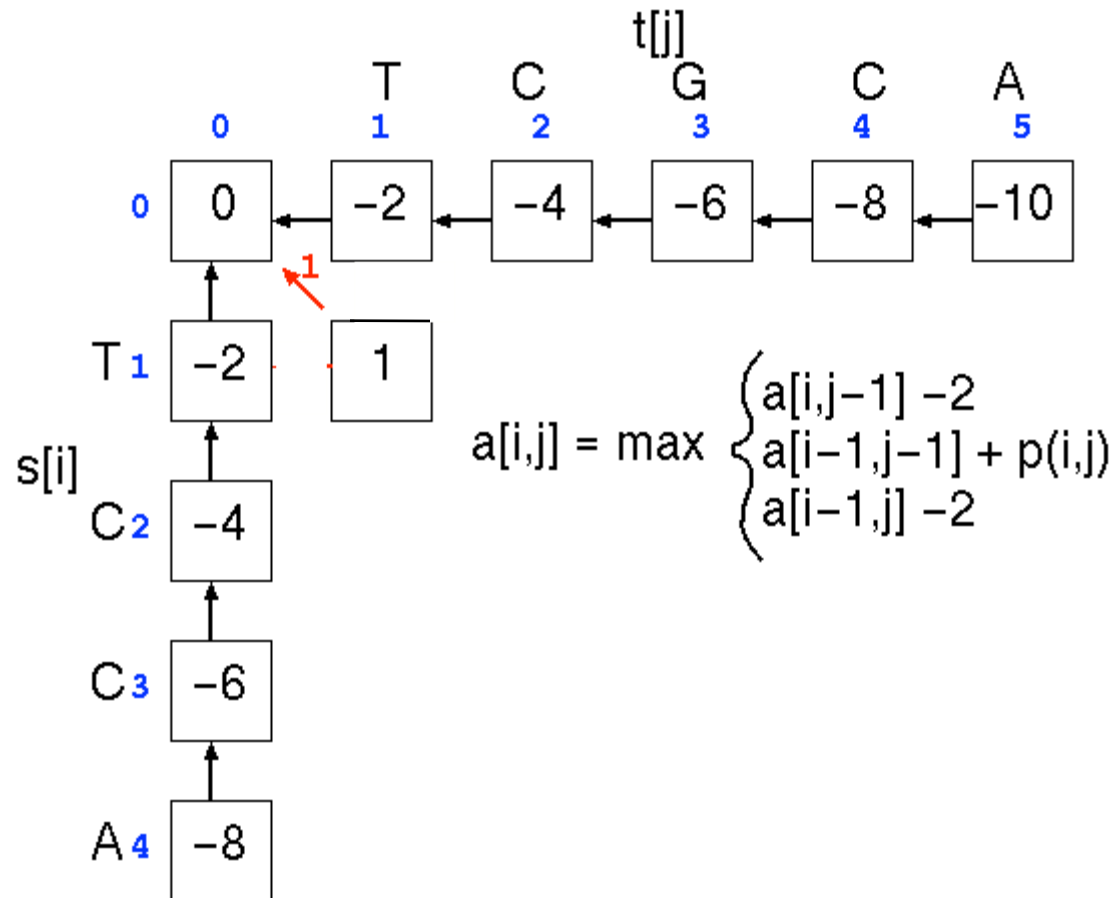
	A	C	G	T
A	1	-1	-1	-1
C	-1	1	-1	-1
G	-1	-1	1	-1
T	-1	-1	-1	1

Gaps: -2

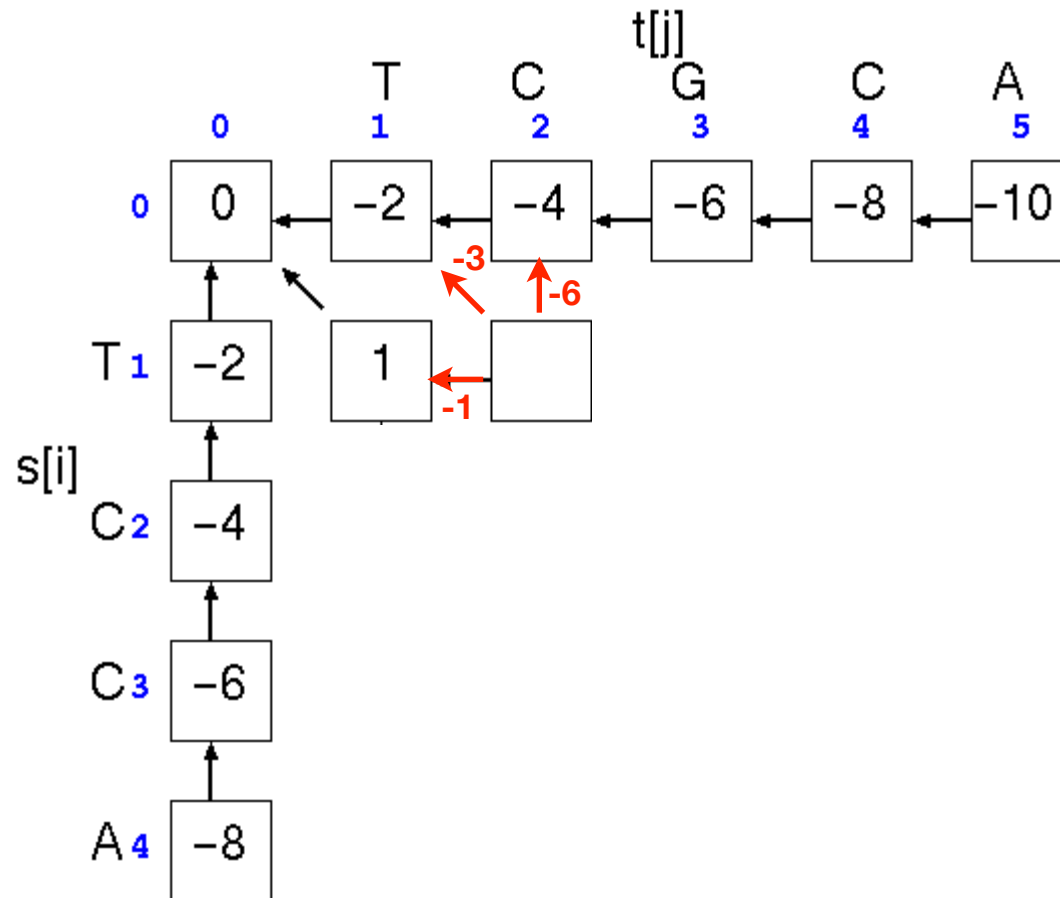
Dynamic programming: example



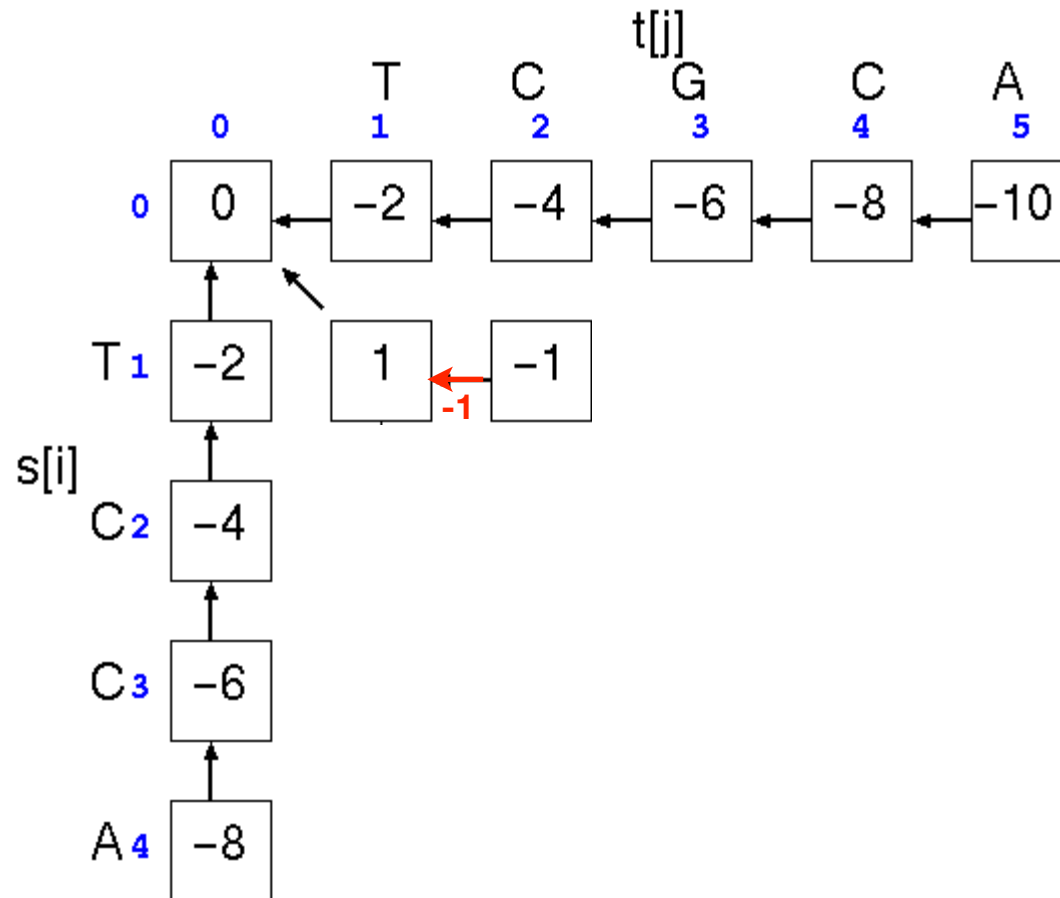
Dynamic programming: example



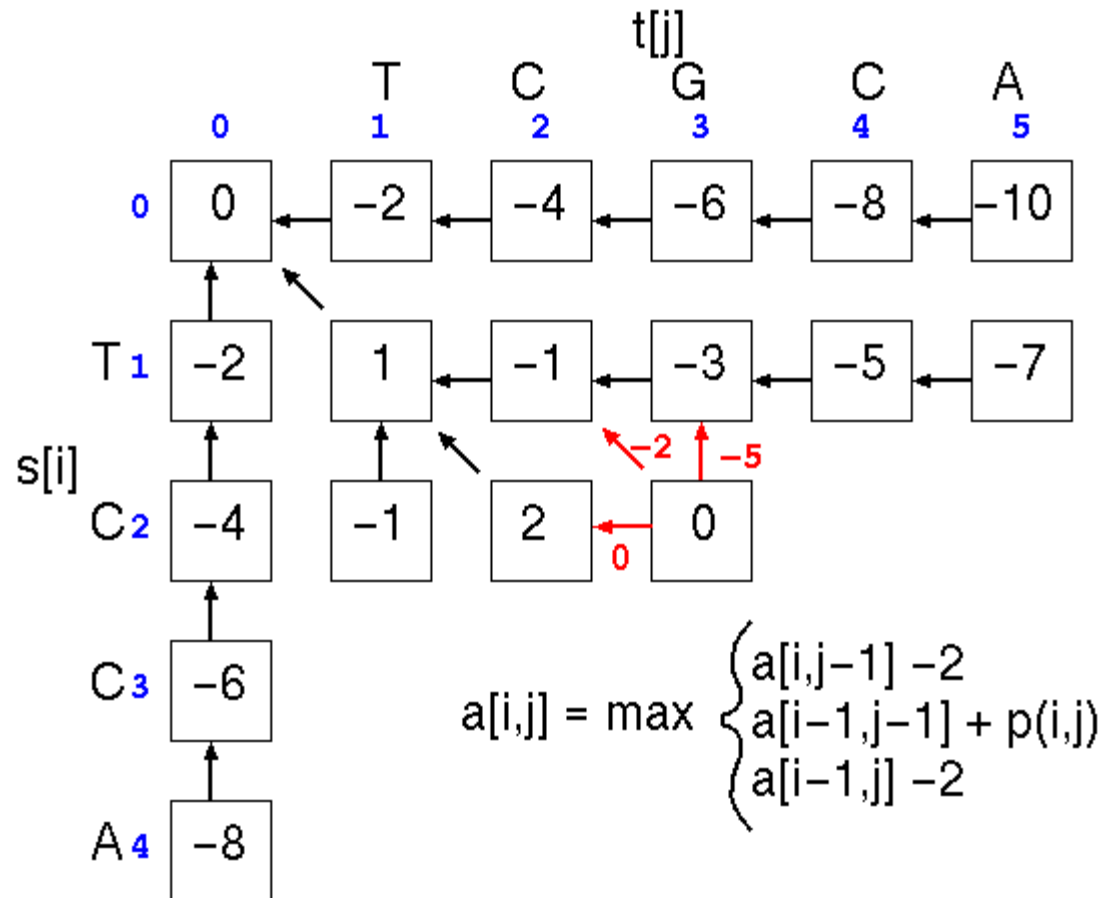
Dynamic programming: example



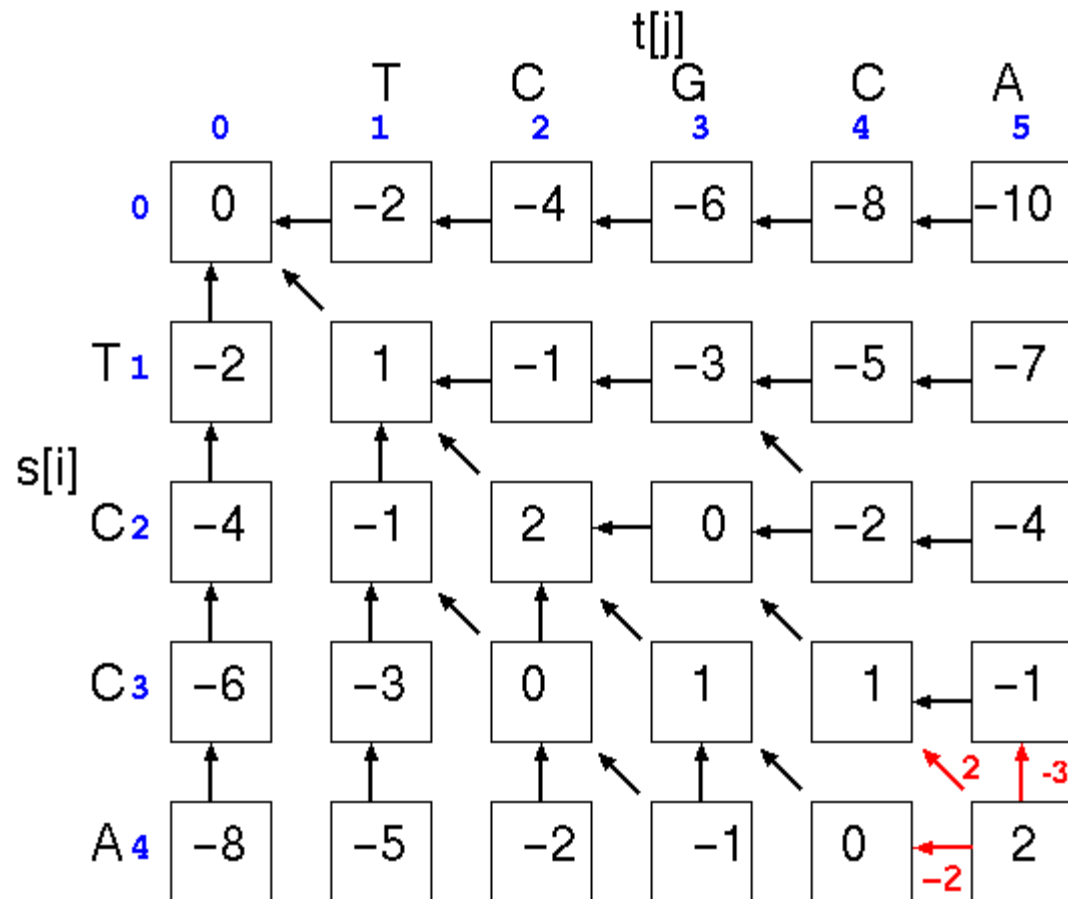
Dynamic programming: example



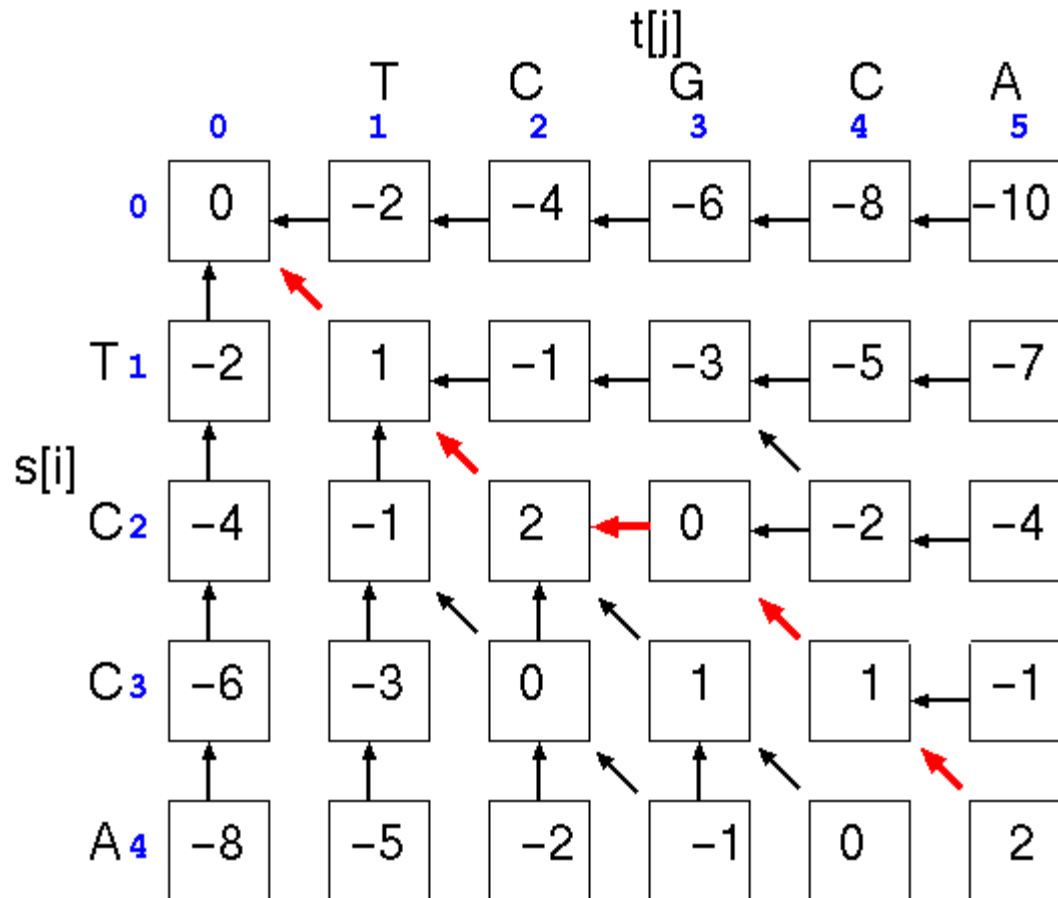
Dynamic programming: example



Dynamic programming: example



Dynamic programming: example



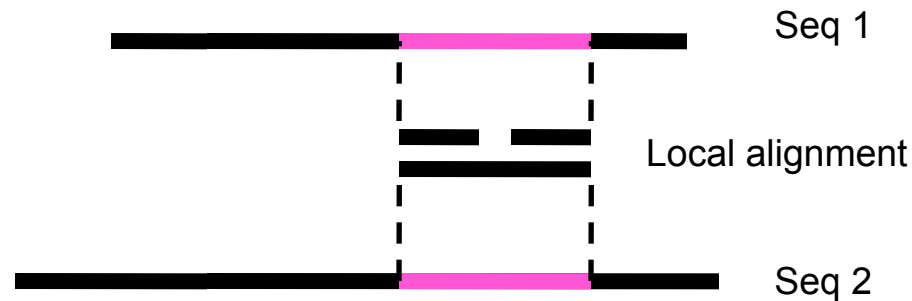
$$\begin{array}{cccccc}
 & T & C & G & C & A \\
 & \vdots & \vdots & & \vdots & \vdots \\
 & T & C & - & C & A \\
 \hline
 & 1 & +1 & -2 & +1 & +1 & = & \underline{2}
 \end{array}$$

Global versus local alignments

Global alignment: align full length of both sequences.
(The “Needleman-Wunsch” algorithm).



Local alignment: find best partial alignment of two sequences
(the “Smith-Waterman” algorithm).



Local alignment overview

- The recursive formula is changed by adding a fourth possibility: zero. This means local alignment scores are never negative.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \\ \text{score}(x-1,y) - \text{gap-penalty} \\ 0 \end{cases}$$

- Trace-back is started at the highest value rather than in lower right corner
- Trace-back is stopped as soon as a zero is encountered

Local alignment: example

		H	E	A	G	A	W	G	H	E	E
	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	5	0	5	0	0	0	0	0
W	0	0	0	0	2	0	20	12	4	0	0
H	0	10	2	0	0	0	12	18	22	14	6
E	0	2	16	8	0	0	4	10	18	28	20
A	0	0	8	21	13	5	0	4	10	20	27
E	0	0	6	13	18	12	4	0	4	16	26

AWGHE

AW-HE

Substitution matrices and sequence similarity

- Substitution matrices come as series of matrices calculated for different degrees of sequence similarity (different evolutionary distances).
- "Hard" matrices are designed for similar sequences
 - Hard matrices are designated by high numbers in the BLOSUM series (e.g., BLOSUM80)
 - Hard matrices yield short, highly conserved alignments
- "Soft" matrices are designed for less similar sequences
 - Soft matrices have low BLOSUM values (45)
 - Soft matrices yield longer, less well conserved alignments

Alignments: things to keep in mind

“Optimal alignment” means “having the highest possible score, given substitution matrix and set of gap penalties”.

This is NOT necessarily the biologically most meaningful alignment.

Specifically, the underlying assumptions are often wrong: substitutions are not equally frequent at all positions, affine gap penalties do not model insertion/deletion well, etc.

Pairwise alignment programs always produce an alignment - even when it does not make sense to align sequences.